# The Earth System Grid Center for Enabling Technologies: Enabling Community Access to Petascale Climate Datasets

**Dean N. Williams, Lawrence Livermore National Laboratory (LLNL)**

**David E. Bernholdt, Oak Ridge National Laboratory (ORNL)**

**Ian T. Foster, Argonne National Laboratory (Argonne)**

**Don E. Middleton, National Center for Atmospheric Research (NCAR)**

## 1. Introduction

Climate research is inherently a multidisciplinary endeavor. As researchers strive to understand the complexity of our climate system they form multi-institutional and multinational teams to tackle "Grand Challenge" problems. These multidisciplinary virtual organizations need a common software infrastructure to access the many large global climate model datasets and tools. It is critical that this infrastructure provide equal access to climate data, supercomputers, simulations, visualization software, whiteboard, and other resources. To this end, we established the Earth System Grid (ESG) Center for Enabling Technologies (ESG-CET) [1], a collaboration of seven U.S. research laboratories (Argonne, LANL, LBNL, LLNL, NCAR, NOAA/PMEL, and ORNL) and a university (USC/ISI) working together to identify and implement key computational and informational technologies for advancing climate change science. Sponsored by the Department of Energy (DOE) Scientific Discovery through Advanced Computing (SciDAC)-2 [2] program, through the Offices of Advanced Scientific Computing Research (OASCR) [3] and the Offices of Biological and Environmental Research (OBER) [4], ESG-CET utilizes and develops computational resources, software, data management, and collaboration technologies to support observational and modeling data archives.

Work on ESG began with the "Prototyping an Earth System Grid" (ESG I) project, initially funded under the DOE Next Generation Internet (NGI) program, with follow-on support from OBER and DOE's Mathematical, Information, and Computational Sciences (MICS) office. In this prototyping project, we developed Data Grid technologies for managing the movement and replication of large datasets, and applied these technologies in a practical setting (an ESG-enabled data browser based on current climate data analysis tools), achieving cross-country transfer rates of more than 500 Mb/s. Having demonstrated the potential for remotely accessing and analyzing climate data located at sites across the U.S., we won the "Hottest Infrastructure" award in the Network Challenge event at the SC'2000 conference.

While the ESG I prototype provided a proof of concept ("Turning Climate Datasets into Community Resources"), the SciDAC Earth System Grid (ESG) II project [5, 6] made this a reality. Our efforts in that project targeted the development of metadata technologies [7] (standard schema, XML metadata extraction based on netCDF, and a Metadata Catalog Service), security technologies [8] (Web-based user registration and authentication, and community authorization), data transport technologies [9, 10] (GridFTP-enabled OPeNDAP-G for high-performance access, robust multiple file transport and integration with mass storage systems, and support for dataset aggregation and subsetting), and web portal technologies to provide interactive access to climate data holdings. At this point, the technology was in place and assembled, and ESG II was poised to make a substantial impact on the climate modelling community.

In 2004, the National Center for Atmospheric Research (NCAR), a premier climate science laboratory, and lead institution for the Community Climate System Model (CCSM) modeling collaboration, began its first publication of climate model data into the ESG system, drawing on simulation data archived at

LANL, LBNL, NCAR, and ORNL. Late that same year, the Program for Climate Model Diagnosis and Intercomparison (PCMDI), an internationally recognized climate data center at LLNL, launched a production service providing access to climate model data germane to the Intergovernmental Panel on Climate Change (IPCC) 4th Assessment Report (AR4) [11]. (Because of international data requirements, restrictions, and timelines, the NCAR and PCMDI ESG data holdings were separated.) ESG has since become a world-renowned leader in developing technologies that provide scientists with virtual access to distributed data and resources.
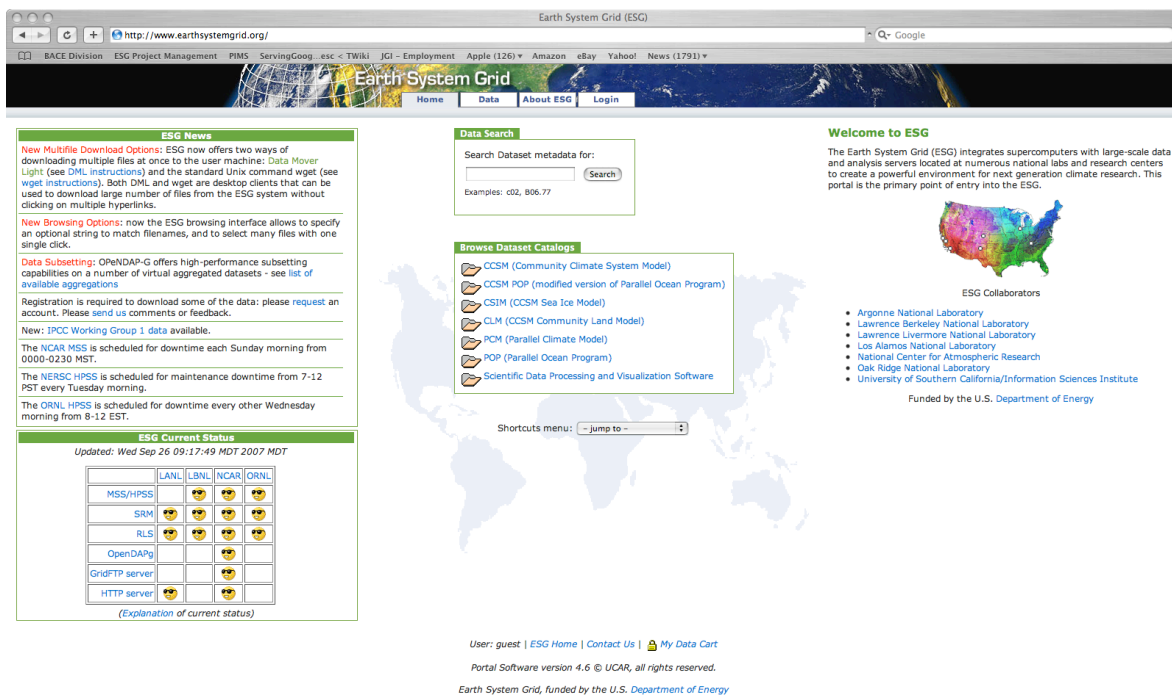
In its first full year of production (late 2005), the two ESG sites provided access to a total of 220 TB of data, served over 3,000 registered users, and delivered over 100 TB of data to users worldwide. Analysis of just one component of ESG data holdings, those relating to the Coupled Model Intercomparison Project phase 3 (CMIP3), resulted in the publication of over 100 peer-reviewed scientific papers.

In 2006 we launched the current phase of the ESG effort, the ESG Center for Enabling Technologies (ESG-CET). The primary goal of this stage of the project is to broaden and generalize the ESG system to support a more broadly distributed, more international, and more diverse collection of archive sites and types of data. An additional goal is to extend the services provided by ESG beyond access to raw data by developing "server-side analysis" capabilities that will allow users to request the output from commonly used analysis and intercomparison procedures. We view such capabilities as essential if we are to enable large communities to make use of petascale data. However, their realization poses significant resource management and security challenges.

## 2. Overview of ESG

ESG is a large, production, distributed system – a Data Grid – with primary access points via three web portals: one for general climate research data, another dedicated to the IPCC activity, and a third for the Community Climate System Model (CCSM) Biogeochemistry (BGC) Working Group which is just going into production at ORNL. The deployment of these three separate portals is driven by international data requirements, restrictions, and timelines. However, they are all based on the same underlying software system. Our goal in ESG-CET is to achieve complete integration of these focused archives, while providing the tailored access and other controls required by the various data owners. In this way, we will provide ESG users with coherent access to ever-growing and increasingly diverse collections of global community climate data.

Users of the ESG portal must first register, at which time they are granted appropriate privileges and access to data collections. The main portal page, shown in Figure 1, provides news, status, and live monitoring of ESG. Once logged in, users may either search or browse ESG catalogs to locate desired datasets, with the option of browsing both collection-level and file/usage-level metadata. Based on this perusal of the catalogs, users may gather a collection of files into a "DataCart" or request an "aggregation," which allows them to request a specific set of variables subject to a spatiotemporal constraint. Selected data may then be downloaded to the user's system, including datasets that are on deep storage at multiple sites behind security firewalls. Group-based authorization mechanisms allow the ESG administrators to control which users can access which data. These capabilities are made possible by a collection of ESG management, data publishing, and large-scale data transport tools.

**Figure 1: ESG Portal**

The ESG system includes a metrics-gathering capability that keeps track of user activity. Interactive displays as well as reports allow us to track what data is downloaded, how often, and by whom. The resulting data has proved invaluable not only for reporting to sponsors and data owners on degree of use (its initial intent) but also as a guide to system development and optimization.
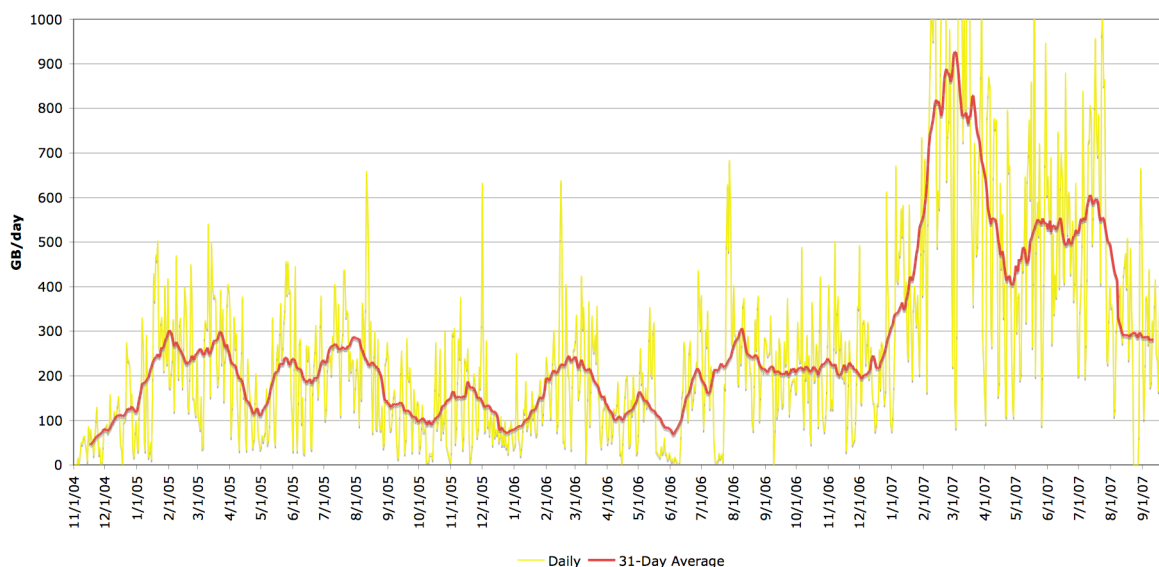
## 3. Overall Impact

ESG has had a significant impact upon the national and international climate community by enabling broad dissemination of important data holdings, including the Community Climate System Model (CCSM) data archive, the Intergovernmental Panel on Climate Change (IPCC) 4th Assessment Report (AR4) data archive, and now the CCSM BGC Carbon-Land Model Intercomparison Project (C-LAMP) [12] data archive. All three archives are well known to the user community and since ESG's official release, the community has downloaded well over 300 TB of data and well over 1 million files, and reported over 300 journal articles [13], all in a short time span.

The ESG team works closely with the CCSM community to publish CCSM model data into the ESG archives. Collaborating with CCSM scientists and data providers, the ESG team developed and utilized Grid technology that interfaces into the ESG metadata database allowing the CCSM community to view and manage all information related to generating, defining, and archiving CCSM model simulation runs. This interface allows scientists to impose selective access control on project runs, to sort information by any type, and to enter data collaboratively. The long-term goal is to tie the metadata ingestion process to the actual CCSM run workflow, so that model simulation metadata can be added automatically into the ESG data holdings.

The ESG user base comprises climate scientists, analysts, educators, governments (both domestic and abroad), private industry, and many others. CCSM data, along with other important datasets accessible via ESG such as those produced by the Parallel Climate Model (PCM) [14] and the Parallel Ocean Program (POP) [15], have been used in numerous scientific papers, impact analyses, urban planning and ecosystem monitoring studies,, education, and other activities. By allowing access, ESG enables scientists, hardware and software engineers, universities and others to

examine and learn how a state-of-the-art climate model works, and to provide suggestions and enhancements for its scientific accuracy, portability, and performance. We even receive occasional queries from the general public, asking how they can use data published in ESG to better understand climate change issues or local impacts.

ESG was thrust into international collaboration when it was asked in late 2003 to support the IPCC/Working Group on Coupled Models (WGCM) need to distribute data to the international climate community. The IPCC, which was jointly established by the World Meteorological Organisation (WMO) and the United Nations Environment Programme, carries out periodic assessments of the science of climate change. Fundamental to this effort is the production, collection and analysis of data from climate model simulations carried out by major international research centers. Analysis of a set of standard climate-change simulations from many modelling centers provides comprehensive understanding of the strengths and weaknesses of climate models, as well as which aspects of the simulation results may be due to characteristics of specific models and which are generally observed across multiple models. The IPCC and WGCM requested that PCMDI at LLNL collect model output data from these IPCC simulations, and distribute these to the community via ESG. Since this effort began, IPCC model runs published to the climate community via the CMIP3 (IPCC AR4) ESG portal total to just over 35 TB (78,158 files), and some 1,400 users have registered to receive IPCC data for analysis. Figure 2 shows the daily download rate over time.



**Figure 2: CMIP3 (IPCC AR4) Download Rates in Gigabytes/day**

New to ESG is the dissemination of C-LAMP [12] biogeochemistry data. This model inter-comparison project has two terrestrial BGC modules linked to the same set of prescribed ocean BGC fluxes, together with the CCSM's interactive atmosphere and interactive land surface modules. The C-LAMP effort involves two separate experiments: one in which atmospheric data comes from observations, the other in which it is calculated by CAM3, the current atmospheric component of the CCSM. The first experiment will determine how well land-air fluxes of $CO_2$ are simulated by the two BGC modules, given the observed climate; the second will determine the effect of the atmospheric model's climate bias (notably in precipitation) on the simulated $CO_2$ fluxes. The C-LAMP experimental output is now being archived and disseminated on an ESG C-LAMP site modelled after the ESG CMIP3 (IPCC AR4). This archive will initially be open only to members of the BGC Working Group, but ultimately the working group will open up the data to any interested researcher.

Knowledge and expertise gained from ESG have helped the climate community plan effective strategies to manage a rapidly growing data environment. Approaches and technologies developed under the ESG project have also impacted data-simulation integration in other disciplines, such as astrophysics, molecular biology, and materials science.

## 4. The Next-Generation ESG

Building upon ESG's success to date, ESG-CET is developing a next-generation environment targeted at enabling flexible, efficient, universal access to yet larger datasets, and to harnessing distributed worldwide resources for the purpose of advancing climate and related impacts research and assessment. In creating this new community infrastructure, ESG-CET will turn even more climate model data into true community resources and place advanced capabilities in the hands of a substantial user base community.

Our high-level goals for this next phase of ESG are driven by scientific objectives relevant to DOE's scientific priorities over the next several years. In brief, they are, firstly, to sustain successful existing ESG services, and secondly, to address scientific needs related to projected future data management and analysis requirements, with a particular focus on:

- Preparing for the CMIP4 IPCC 5th Assessment Report (AR5) in 2010.
- Publishing and enabling processing of the massive data produced by the Climate Science Computational End Station (CCES) at ORNL's NCCS/LCF.
- Supporting a wide-range of climate model evaluation activities aimed at improving climate change research.

To support this effort, we will broaden ESG to support multiple types of model and observational data, provide more powerful (client-side) ESG access and analysis services, enhance interoperability between common climate analysis tools and ESG, and enable end-to-end simulation and analysis workflow. Figure 3 depicts the scientific data management and analysis requirements in relationship to the ESG development timeframe. We specifically note that a distributed testbed for CMIP4 (IPCC AR5) must be in place by early 2009.
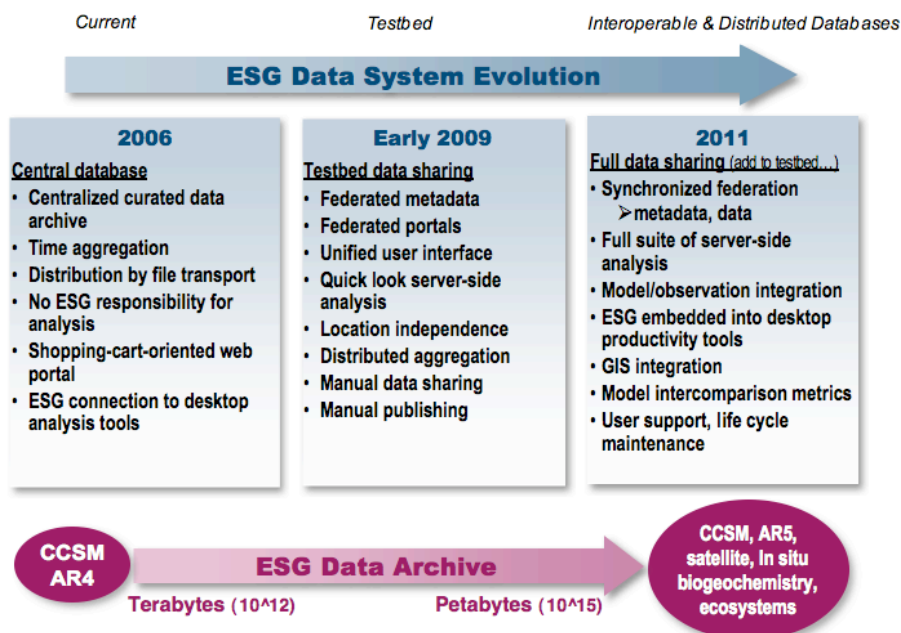


**Figure 3: Evolving ESG to the Petascale: High-level ESG-CET Roadmap**

The ESG-CET architecture must be generalized to enable a larger number of sites with more diverse capabilities to selectively federate, cooperate, or operate in a standalone fashion as individual sites desire. The architecture must support a variety of user access mechanisms, including multiple portals and service- or API-based access, and data delivery mechanisms. This architecture must also be robust in the face of system and network failures at the participating sites.

To address these concerns, we designed the federated ESG-CET architecture (see Figures 4 and 5) to provide interoperability and enhanced functionality to users, and are now implementing the new design through a combination of evolution of existing software, development of new tools, and integration with third-party software. The much wider deployment anticipated for the next generation system means that software deployability and maintainability are vital considerations in determining the most effective implementation.



**Figure 4: Future ESG-CET Architecture**



**Figure 5: The ESG-CET Federated System**

## 5. Conclusion

ESG has made significant progress towards the definition of the federated metadata, security, and data services required to enabled distributed access to, and analysis of, large quantities of climate simulation data. The current production-level ESG system primarily addresses the tasks of publishing and cataloging terabytes of climate model data for a diverse set of registered users. We are now working to take ESG-CET to the next level of distributed environments with an even greater emphasis on federation and server-side capabilities. ESG-CET will build upon the current ESG system and target flexibility, efficiency, and more universal access while expanding to serve much larger archives (petabytes), as required for CMIP4 (IPCC AR5), CCSM, and CCES. To this end, ESG-CET is working with disparate technologies and partnering with national and international leaders in the computer and climate communities to build a robust data and analysis distributed infrastructure in support of advancing climate change research.

## 7. References

1. Earth System Grid (ESG) - Turning Climate Model Datasets into Community Resources. http://www.earthsystemgrid.org/, 2004 – 2007.
2. Scientific Discovery through Advanced Computing (SciDAC). http://www.scidac.gov/, 2007.
3. Office of Advanced Scientific Computing Research (OASCR). http://www.science.doe.gov/ascr/, 2007.
4. Offices of Biological and Environmental Research (OBER). http://www.science.doe.gov/ober/, 2007.
5. Bernholdt, D., Bharathi, S., Brown, D., Chanchio, K., Chen, M., Chervenak, A., Cinquini, L., Drach, B., Foster, I., Fox, P., Garcia, J., Kesselman, C., Markel, R., Middleton, D., Nefedova, V., Pouchard, L., Shoshani, A., Sim, A., Strand, G. and Williams, D. The Earth System Grid: Supporting the Next Generation of Climate Modeling Research. *Proceedings of the IEEE*, *93* (3). 485-495. 2005.
6. Foster, I., Alpert, E., Chervenak, A., Drach, B., Kesselman, C., Nefedova, V., Middleton, D., Shoshani, A., Sims, A. and Williams, D., The Earth System Grid: Turning Climate

Datasets Into Community Resources. in *82nd Annual American Meteorological Society Meeting*, (Orlando, FL., 2002).

7. Eaton, B., Gregory, J., Drach, B., Taylor, K. and Hankin, S. NetCDF Climate and Forecast Metadata Conventions. http://cf-pcmdi.llnl.gov, 2007.

8. Welch, V., Siebenlist, F., Foster, I., Bresnahan, J., Czajkowski, K., Gawor, J., Kesselman, C., Meder, S., Pearlman, L. and Tuecke, S., Security for Grid Services. *12th IEEE International Symposium on High Performance Distributed Computing*. 2003.

9. Fox, P., Garcia, J. and West, P. OPeNDAP for the Earth System Grid. *Data Science Journal*. 2007.

10. Allcock, B., Bresnahan, J., Kettimuthu, R., Link, M., Dumitrescu, C., Raicu, I. and Foster, I., The Globus Striped GridFTP Framework and Server. *Supercomputing 2005 (SC '05) conference proceedings*. 2005.

11. Intergovernmental Panel on Climate Change (IPCC) 4th Assessment Report. http://www.ipcc.ch/activity/ar.htm, 2007.

12. CCSM Carbon LAnd Model intercomparison Project (C-LAMP). http://www.climatemodeling.org/c-lamp/, 2007

13. World Climate Research Program (WCRP) CMIP3 (IPCC AR4) Subproject Publications. http://www-pcmdi.llnl.gov/ipcc/subproject_publications.php, 2007.

14. The Parallel Climate Model. http://www.cgd.ucar.edu/pcm/. 2007

15. The Parallel Ocean Program (POP) ocean circulation model. http://climate.lanl.gov/Models/POP/. 2007.